
[Re] Subspace Attack: Exploiting Promising Subspaces for Query-Efficient Black-box Attacks

Edoardo Debenedetti
École polytechnique fédérale de Lausanne
edoardo.debenedetti@epfl.ch

Dan Shaked Renous
Weizman Institute of Science
dan.shakedrenous@weizmann.ac.il

Faezeh Nassajian
École polytechnique fédérale de Lausanne
faezeh.nassajian@epfl.ch

Abstract

As part of the NeurIPS 2019 Reproducibility Challenge and EPFL’s CS-443 Machine Learning course, we chose to attempt reproduce the attack algorithm proposed in “Subspace Attack: Exploiting Promising Subspaces for Query-Efficient Black-box Attacks”. Our reported results are better than the original paper in terms of the median number of queries per attack, but worse in terms of failure rate. A concise assessment of our implementation is also included.

1 Introduction

As the use of Machine Learning in services and applications is increased, it becomes important to assess its reliability and security. In the last few years, several flaws have been found in state-of-the-art Machine Learning methods and algorithms. One of the most studied flaws is the fact that many models can be attacked using the so-called adversarial examples, inputs that are minimally perturbed in order to be misclassified by the victim model.

In this paper we attempt to reproduce the results obtained in the paper “Subspace Attack: Exploiting Promising Subspaces for Query-Efficient Black-box Attacks” by Yan et al. [1], published among the proceedings of NeurIPS 2019. In their paper, the authors present a new kind of black-box attack, that, for the first time, ensembles attacks’ transferability and gradient estimation for Projected Gradient Descent.

Our report has been written for the NeurIPS 2019 Reproducibility Challenge¹ and EPFL’s² CS-443 Machine Learning course. It consists of a background about adversarial attacks (Section 2), our methodology (Section 3) and a concise reproducibility section (Section 4). These are followed by our results, discussion, and conclusions (Sections 5 to 7, respectively).

2 Background

2.1 Adversarial Examples

Adversarial examples are inputs fed to a machine learning model which mislead the model to make an incorrect prediction. Several machine learning architectures have been proven to be vulnerable in adversarial settings, where the adversarial examples can be generated with several attacking

¹<https://reproducibility-challenge.github.io/neurips2019/>

²École polytechnique fédérale de Lausanne

techniques. Thus, the aim of such attack is to perturb an input $\mathbf{x} \in \mathbb{R}^n$ and to trick a victim model $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$, such that it gives a wrong prediction, $\arg \max_i f(\mathbf{x})_i \neq y$ (where y is the true label). Such general definition is called *untargeted*. A more specific type of attack is a *targeted* one, which aims towards a misleading classification of a specific label $y' \neq y$, $\arg \max_i f(\mathbf{x})_i = y'$.

In general, adversarial attacks can be classified in two main types: white-box and black-box. While only the confidence score from a target victim model is accessible to the latter, the former has full access to the network parameters of it. This allows an efficient use of various gradient-based methods [2] to generate adversarial examples. This emphasises the attacker’s preference for the white-box models. However, it is not a realistic setting.

Let us now consider the more realistic setting in which the attacker has no access to the trained victim model parameters. In this case, it is not possible to compute the gradient of the model with respect to the perturbed input. Nonetheless, several methods have been proposed to overcome this issue. Some of these methods rely on the so-called adversarial example transferability [3], in which a surrogate model is trained with data samples labeled by the victim model, and is then used to generate adversarial examples that work on the victim model, too. The other kind of methods make use of first-order optimization methods (e.g. finite differences, as proposed in [4]) by querying the victim model to estimate the gradient of the victim model to run Projected Gradient Descent (PGD). A drawback of these methods is the need of a large number of queries for a good estimation of the victim model’s gradient. This motivates the attempts for reducing the required number of queries for a successful attack. The next subsection describes the solution proposed in [1], which harnesses the features of both zeroth-order optimization and transferability.

2.2 Subspace Attack Algorithm

In what follows, we present the essence of the subspace attack. This will include introducing the principle of the *Bandit attack* [5], on which the subspace attack heavily relies on.

Given a classification loss function $\mathcal{L}(x, y)$, where x is some input and y its corresponding label. We can formulate the adversarial attack problem as follows:

$$x' = \arg \max_{x': \|x' - x\|_p \leq \epsilon_p} \mathcal{L}(x', y) \quad (1)$$

Let $g^* = \nabla_x \mathcal{L}(x, y)$ be the gradient of \mathcal{L} at (x, y) , Then the goal of the gradient estimation problem is to find a unit vector \hat{g} that maximizes the inner product $\mathbb{E} [\hat{g}^T g^*]^3$.

In [5], the authors have proven that the gradient estimation problem can be considered as a bandit optimization problem. Bandit optimization is a tool used in online convex optimization, in which there is an agent playing a game that includes a sequence of rounds. During each round t , the agent must select an action and incurs in a loss \mathcal{L}_t , whose expectation across all the rounds should be minimized. The main novelty introduced by [5] is the fact that the estimation is improved by accumulating prior information about gradients in a latent vector \mathbf{g}_t , that is updated each round with an estimation of the gradient of the victim model, performed via finite differences. In [5], each new basis \mathbf{u}_t used for the gradient estimation is randomly sampled from a Gaussian distribution.

In the subspace attack proposal, the authors present a novel way to accumulate prior information about the gradient. Instead of using the full basis of the sample’s space to estimate the true gradient, a subspace of directions is chosen by using the gradients of reference models. These models are pre-trained on the same dataset as the victim and their parameters are fully known. Each round of the bandit optimization uses the gradient of a randomly chosen reference model, with respect to the adversarial input, as a basis for the gradient estimation. Moreover, [1] proposes a way to apply different ratios of drop-out to the enrich the set of reference models.

3 Methodology

In this work, we implemented the subspace attack using PyTorch [6], following the algorithm specified in [1], without looking at the source code of the attack in the original paper. We chose to evaluate the reproducibility of the algorithm by attempting to replicate untargeted subspace attacks

³The expectation here is taken over the randomness of the estimation algorithm

Algorithm 1: Subspace Attack, as described in [1].

Input: A benign example $\mathbf{x} \in \mathbb{R}^n$, its label y , a set of m reference models $\{f_0, \dots, f_{m-1}\}$, a chosen attack objective function $\mathcal{L}(\cdot, \cdot)$, and a victim model from which the output of f can be inferred.

Output: An adversarial example \mathbf{x}_{adv} that fulfills $\|\mathbf{x}_{adv} - \mathbf{x}\|_\infty \leq \epsilon$.

- 1: Initialize the adversarial example to be crafted $\mathbf{x}_{adv} \leftarrow \mathbf{x}$.
 - 2: Initialize the gradient to be estimated $\mathbf{g} \leftarrow \mathbf{0}$.
 - 3: Initialize the drop-out/layer ratio p .
 - 4: **while** not successful **do**
 - 5: Choose a reference model whose index is i , uniformly at random.
 - 6: Calculate prior gradient with drop-out/layer ratio p as $\mathbf{u} \leftarrow \frac{\partial \mathcal{L}(f_i(\mathbf{x}_{adv}, p), y)}{\partial \mathbf{x}_{adv}}$
 - 7: $\mathbf{g}_+ \leftarrow \mathbf{g} + \tau \mathbf{u}$, $\mathbf{g}_- \leftarrow \mathbf{g} - \tau \mathbf{u}$
 - 8: $\mathbf{g}'_+ \leftarrow \mathbf{g}_+ / \|\mathbf{g}_+\|_2$, $\mathbf{g}'_- \leftarrow \mathbf{g}_- / \|\mathbf{g}_-\|_2$
 - 9: $\Delta_t \leftarrow \frac{\mathcal{L}(f(\mathbf{x}_{adv} + \delta \mathbf{g}'_+), y) - \mathcal{L}(f(\mathbf{x}_{adv} + \delta \mathbf{g}'_-), y)}{\tau \delta} \mathbf{u}$
 - 10: $\mathbf{g} \leftarrow \mathbf{g} + \eta \mathbf{g} \Delta_t$
 - 11: $\mathbf{x}_{adv} \leftarrow \mathbf{x}_{adv} + \eta \cdot \text{sign}(\mathbf{g})$
 - 12: $\mathbf{x}_{adv} \leftarrow \text{Clip}(\mathbf{x}_{adv}, \mathbf{x} - \epsilon, \mathbf{x} + \epsilon)$
 - 13: $\mathbf{x}_{adv} \leftarrow \text{Clip}(\mathbf{x}_{adv}, 0, 1)$
 - 14: Update the drop-out/layer ratio p , following the original paper’s policy.
 - 15: **end while**
 - 16: **return** \mathbf{x}_{adv}
-

on the GDAS [7] and WRN [8] models trained on the CIFAR-10 dataset [9]. We use as reference models VGG-11/13/16/19 [10] and AlexNet [11]. In order to verify that we use the same settings – hyper-parameters and pre-trained models – as the original subspace experiments, we checked their log files and loaded the pre-trained models from their source-code repository. The latter was also required since the available models by PyTorch are pre-trained on ImageNet, while our experiment uses images from CIFAR-10.

Aiming to further assess the performance of our algorithm implementation, we added options to keep track of some values encountered during the attack. In addition to the required number of queries for a successful attack of each image classification, we recorded in each iteration the cosine similarity between the estimated gradient and the true gradient (as done in [5]), and the norms of these gradients.

Using the principle of evaluation of [1], the evaluation of an attack’s performance is done with respect to its failure rate and the number of queries. Since the distribution of the latter varies, we use the mean and the median number of queries.

4 Reproducibility

We first try to implement Algorithm 1, and run the experiment attacking GDAS, using VGGNets and AlexNet as reference models, using the hyper-parameters listed in [1] and [5]. We then expand our implementation to be used to attack WRN and Pyramidnet, and look for better hyper-parameters.

4.1 Machine Setup, experiment duration and budget

We run our experiment both on Google Cloud Platform and on Code Ocean [12]. The Ubuntu Virtual Machine on Google Cloud Platform has an Nvidia Tesla T4 GPU, 52 GB memory and 8 vCPUs. A set of 1000 attacks against GDAS using VGGs and AlexNet as reference models took us about 7h45m to run, with a total of \$7.75 spent. On Code Ocean, the Virtual Machine we used runs Ubuntu on a 4-cores CPU with 60 GB of memory and an Nvidia Tesla K80 GPU. As Code Ocean is a sponsor of the Reproducibility Challenge, we have been provided with free compute time to run the experiments. A set of 1000 attacks against GDAS using VGGs and AlexNet as reference models took us about 9 hours to run.

Even though our reported results are produced using GPUs, our implementation can run on a CPU as well. This was tested on a laptop with Ubuntu, an Intel Core i7-7500U and 8GB memory. The computation time takes about 6 times longer than a GPU (in terms of iterations per second). In addition, the results obtained could be different, due to differences in floating point precision and implementations of low-level operations⁴. Such results are outside of the scope of this challenge, and thus not reported.

4.2 Algorithm implementation

The source code of our implementation is available online⁵. As mentioned in section 3, our implementation is done using the PyTorch framework. We use the pre-trained reference and victim models⁶ provided with the original paper’s code repository, along with their corresponding Python classes.

The CIFAR-10 dataset is loaded using PyTorch’s `torchvision.datasets`, and is iterated using a `DataLoader`. The dataset is shuffled, but the seed can be fixed, for comparable results. In all our experiments, for the sake of reproducibility, both PyTorch and NumPy seeds are set to 0, and the following values are set: `torch.backends.cudnn.deterministic = True` and `torch.backends.cudnn.benchmark = False`.

Regarding the implementation of the attack itself, we leveraged PyTorch autograd capabilities to compute the gradient of the reference models. As there was no indication about the loss function to be used, we implemented the algorithm to give the user the possibility to use both the Cross Entropy and Negative Log Likelihood functions. However, as the initial results we obtained were extremely similar, we ended up using Cross Entropy for all the reported results.

Two clipping procedures of the adversarial example are presented in lines 12 and 13 of Algorithm 1. While the former was implemented using PyTorch’s functions `torch.min` and `torch.max`, in order to keep \mathbf{x}_{adv} in the region included in $[\mathbf{x} - \epsilon, \mathbf{x} + \epsilon]$ with ℓ_∞ -norm, the latter was implemented using PyTorch’s `torch.clamp`, with 0 and 1 as arguments.

4.3 Hyper-parameters and experiments settings

The hyper-parameters used in a subspace attack⁷ are specified in Table 1. In addition to these, there’s the dropout/layer ratio of the reference models (p).

In [1], it is claimed that the used hyper-parameters were the same as those used in [5]. However, we have found out that the notation was different and a bit confusing. We reconstructed a translation in Table 1 by cross-referencing the notation used in Algorithm 1 of [1] and in Algorithms 1, 2 and 3 of [5]. The reader should note that we have listed two letters as translations of the δ used in [1] – ϵ and η . This is due to the fact that, in Algorithm 2 of [5], an ϵ is used as finite difference probe, but in the Table 3 of Appendix C of the same paper, where hyper-parameters values are listed, the finite difference probe is listed as “ η (Finite difference probe)”.

Table 1: Translation between hyper-parameters in Subspace attack [1] and the bandit attack [5].

	Subspace Attack [1]	Bandit Attack [5]	Value [1]
Bandit exploration	τ	δ	1.0
Finite difference probe	δ	ϵ/η	0.1
Image ℓ_p learning rate	η	h	1/255
OCO learning rate	η_g	η	100

To choose the hyper-parameters’ values, we have started from those stated in [1], which should be the same used in [5], excluding ϵ and η . However, as discussed in Section 6, we obtained results which were far worse than those obtained in [1]. After exploring the published logs of the experiment⁸ we

⁴<https://pytorch.org/docs/stable/notes/randomness.html>

⁵<https://github.com/epfl-ml-reproducers/subspace-attack-reproduction>

⁶<https://go.epfl.ch/subspace-pretrained>

⁷For hyper-parameters, we use the same notation as [1]

⁸<https://go.epfl.ch/subspace-original-logs>

found out that in the original paper, a value of $\eta_g = 0.1$, instead of $\eta_g = 100$ was used. Thus the set of hyper-parameters which yield the best results is that presented in 1 with $\eta_g = 0.1$.

The dropout was applied to the reference models, in the following procedure: each experiment starts with a 0.05 dropout ratio which is increased every iteration by 0.01, until a maximum ratio of 0.5 is reached.

Finally, using the above settings, each experiment included untargeted attacks on the classification of 1,000 images using a limit of 10,000 queries per image – after which the attack is considered a failure. Moreover, we set a maximum perturbation of $\epsilon = 8/255$, defined using ℓ_∞ -norm.

5 Results

In this section, we present the assessment of our implementation as mentioned in Section 3, followed by a comparison of the effectiveness of our subspace attacks to those reported in the original paper.

In order to evaluate the performance of the algorithm and its implementation, we saved the progression of the loss values along with the estimated and true gradients information of a set of 49 attacks on GDAS. We required each attack to complete 5,000 gradient estimation iterations, i.e. 10,000 queries. Figures 1a and 1b presents the mean value, along the gradient estimation process, of the cross entropy loss of the victim model and the cosine similarity between the estimated and the true gradient, respectively. Having only a single failed attack in this run, its corresponding curve presents the exact value, while the rest of the curves present value averaged on the 48 successful attacks and 49 attacks overall. The fluctuations of the curves are a feature of the stochastic nature of the algorithm.

In addition, we scanned the values of the η_g hyper-parameter. As presented in Table 2, it seems that $\eta_g = 0.1$ yields the best results in terms of mean and median for both victim models. Our best results of the different attack experiments are compared with the two baselines [1, 5] in Table 3. There, it is demonstrated that while our implementation managed to get similar mean and better median numbers of queries as the subspace attack paper, it leads to higher failure rates with respect to the baselines.

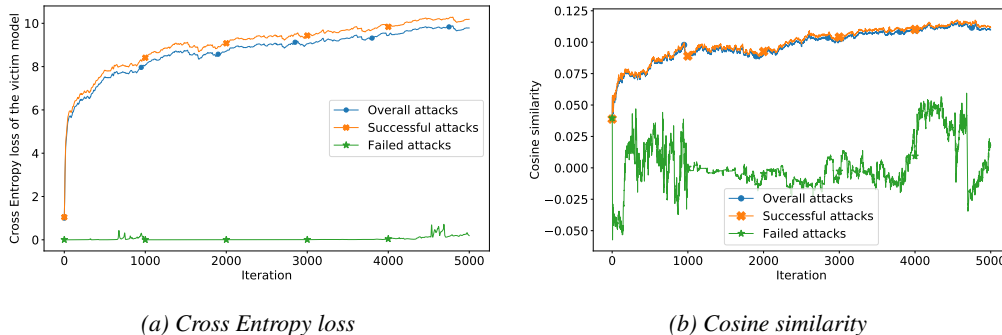


Figure 1: The mean value of the Cross Entropy loss of the true victim model (a) of cosine similarity between the estimated and the true gradient (b) as a function of the estimation iteration in our Subspace Attack implementation. The results of the successful attacks, the failed attacks, and the overall were averaged over 48, 2, and 50 respectively.

Table 2: Comparison between the results obtained with different η_g .

Victim	η_g	Mean Queries	Median Queries	Failure Rate
WRN	0.01	325	16	3.8%
	0.1	330	14	3.7%
GDAS	0.001	263	18	6.3%
	0.01	321	20	5.9%
	0.1	282	18	6.7%
	100	354	16	12.7%

Table 3: Comparison between the results obtained by our implementation and those of the original paper [1] and of [5]. For the results obtained with the other implementations we report the results presented in [1]. In η_g for the results obtained by the original paper we state 0.1* as this is the value we found in their logs, and not the one stated in their paper.

Victim Model	Method	Ref. Models	η_g	Mean Queries	Median Queries	Failure Rate
WRN	Bandit-TD [5]	-	100	713	266	1.2%
	Subspace Attack [1]	AlexNet+VGGNets	0.1*	392	60	0.3%
	Subspace Attack (Ours)	AlexNet+VGGNets	0.1	330	14	3.7%
GDAS	Bandit-TD [5]	-	100	373	128	0.0%
	Subspace Attack [1]	AlexNet+VGGNets	0.1*	250	58	0.0%
	Subspace Attack (Ours)	AlexNet+VGGNets	0.1	282	18	6.7%

6 Discussion

The trends of the mean values of the loss and the cosine similarity (Figures 1a and 1b) confirms that our implementation is indeed efficiently attacking the victim model, as the loss is increased, and the directional agreement between the estimated and the true gradient improves. Furthermore, the obtained results of cosine similarity exceed the ones presented in [5], since our successful attacks demonstrate twice the cosine similarity than those reported in the reference, in the first 5,000 iterations.

Keeping the guiding rule of [1], we use the hyper-parameters of [5], excluding ϵ and η . As we are performing attacks that make use of ℓ_∞ -norm, we should use $\eta_g = 100$. This yields worse results (c.f.r. table 2, and indeed the logs of the subspace paper mention the use of $\eta_g = 0.1$. This also indicates that the performance of our implementation is close to the original one).

Comparing our results to those of the original paper, we notice that while we managed to reduce the computational cost of the attacks, we couldn't reach the reported failure rates of either the original paper or the Bandit attack. Even though we verified that we followed every step of the algorithm, controlled for difference in input models and hyper-parameters, we couldn't improve the originally produced failure rate.

7 Conclusion

In this work, we re-implemented the algorithm of the subspace black-box adversarial attack presented by Yan et al. [1], and introduced methods for the evaluation of its performance. However, we didn't manage to match our effectiveness to that of the original paper. In terms of the reproducibility of [1], it is advisable for the authors to include a clearer statement regarding the used set of hyper-parameters. In addition, when relying on settings from another work, a corresponding mapping of the notation changes can be highly beneficial.

Acknowledgments

The authors would like to thank the NeurIPS Reproducibility Challenge sponsor Code Ocean [12], for providing them with additional compute time for free.

References

- [1] Y. Guo, Z. Yan, and C. Zhang, “Subspace attack: Exploiting promising subspaces for query-efficient black-box attacks,” in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 3820–3829. [Online]. Available: <http://papers.nips.cc/paper/8638-subspace-attack-exploiting-promising-subspaces-for-query-efficient-black-box-attacks.pdf>
- [2] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” *arXiv preprint arXiv:1312.6199*, 2013.
- [3] N. Papernot, P. D. McDaniel, and I. J. Goodfellow, “Transferability in machine learning: from phenomena to black-box attacks using adversarial samples,” *CoRR*, vol. abs/1605.07277, 2016. [Online]. Available: <http://arxiv.org/abs/1605.07277>
- [4] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, “Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models,” in *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, ser. AISC ’17. New York, NY, USA: ACM, 2017, pp. 15–26. [Online]. Available: <http://doi.acm.org/10.1145/3128572.3140448>
- [5] A. Ilyas, L. Engstrom, and A. Madry, “Prior convictions: Black-box adversarial attacks with bandits and priors,” *ICLR 2019*, 2018. [Online]. Available: <https://arxiv.org/abs/1807.07978>
- [6] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- [7] X. Dong and Y. Yang, “Searching for a robust neural architecture in four gpu hours,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1761–1770.
- [8] S. Zagoruyko and N. Komodakis, “Wide residual networks,” *arXiv preprint arXiv:1605.07146*, 2016.
- [9] A. Krizhevsky, G. Hinton *et al.*, “Learning multiple layers of features from tiny images,” University of Toronto, Tech. Rep., 2009.
- [10] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [11] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [12] A. Clyburne-Sherin, X. Fei, and S. A. Green, “Computational reproducibility via containers in social psychology,” *Meta-Psychology*, vol. 3, 2019.